

Lasso estimation of an interval-valued multiple regression model

Marta García Bárzana¹, Ana Colubi¹, and Erricos J. Kontoghiorghes²

¹ Department of Statistics. University of Oviedo
C/ Calvo Sotelo s/n, Oviedo, 33007, Spain

² Department of Commerce, Finance and Shipping
Cyprus University of Technology. P.O. Box 50329, CY-3603 Limassol, Cyprus
{garciabmarta@uniovi.es,colubi@uniovi.es}
erricos@cut.ac.cy

Abstract. A multiple interval-valued linear regression model considering all the cross-relationships between the midpoints and spreads of the intervals has been introduced recently. A least-squares estimation of the regression parameters has been carried out by transforming a quadratic optimization problem with inequality constraints into a linear complementary problem and using Lemke's algorithm to solve it. Due to the irrelevance of certain cross-relationships, an alternative estimation process, the Lasso, is developed. A comparative study showing the differences between the proposed estimators is provided.

Keywords: Multiple regression, Lasso estimation, interval data

1 Introduction

Intervals represent a powerful tool to capture the imprecision of certain characteristics that cannot be fully described with a real number. For example, the measures provided by instruments which have some errors in their measurements [1]. Moreover, intervals also model some features which are inherently interval-valued. For instance, the range of variation of the blood pressure of a patient along a day [2] or the tidal fluctuation [9].

The statistical study of regression models for interval data has been extensively addressed lately in the literature [2–5, 7], deriving into several alternatives to tackle this problem. On one hand, the estimators proposed in [4, 7] account the non-negativity constraints satisfied by the spread variables, but do not assure the existence of the residuals. Hence, they can lead to ill-defined estimated models. On the other hand, the models proposed in [2, 3, 5] are formalized according to the natural interval arithmetic and their estimators lead to models that are always well-defined over the sample range.

The multiple linear regression model [3] considered belongs to the latter approach and its main advantage is the flexibility derived from its way to split the regressors, allowing us to account for all the cross-relationships between the centers and the radii of the interval-valued variables. Nevertheless, this fact

entails an increase in the number of regression parameters and thus, a Lasso estimation is considered in order to shrink some of these coefficients towards zero. The Lasso estimation of an interval-valued regression model has been previously addressed in [4], but being this a more restrictive model formalized in the first framework.

The paper is organized as follows. Section 2 presents some preliminary concepts about the interval framework and Section 3 contains the formalization of the model. The Least-Squares and Lasso estimations of the proposed model are developed in subsections 3.1 and 3.2. Section 4 briefly describes the Lasso model proposed by Giordani [4]. The empirical performance of the estimators proposed in Sections 3 and 4 is compared in Section 5 by means of a illustrative real-life example. Section 6 finishes with some conclusions.

2 Preliminaries

Interval data are defined as elements belonging to the space $\mathcal{K}_c(\mathbb{R}) = \{[a_1, a_2] : a_1, a_2 \in \mathbb{R}, a_1 \leq a_2\}$. Given an interval $A \in \mathcal{K}_c(\mathbb{R})$, it can be parametrized in terms of its center or *midpoint*, $\text{mid } A = (\sup A + \inf A)/2$, and its radius or *spread*, $\text{spr } A = (\sup A - \inf A)/2$. Nonetheless, intervals can alternatively be expressed by means of the so-called canonical decomposition [2] given by $A = \text{mid } A[1 \pm 0] + \text{spr } A[0 \pm 1]$. This decomposition allows us to consider separately the *mid* and *spr* components of A , which will lead into a more flexible model. The interval arithmetic on $\mathcal{K}_c(\mathbb{R})$ consists of the Minkowski addition and the product by scalars defined as follows by the jointly expression: $A + \delta B = [(\text{mid } A + \delta \text{mid } B) \pm (\text{spr } A + |\delta| \text{spr } B)]$ for any $A, B \in \mathcal{K}_c(\mathbb{R})$ and $\delta \in \mathbb{R}$.

The space $(\mathcal{K}_c(\mathbb{R}), +, \cdot)$ is not linear but semilinear, as the existence of symmetric element with respect to the addition is not guaranteed in general. An additional operation is introduced, the so-called Hukuhara difference between the intervals A and B . The difference C is defined as $C = A -_H B \in \mathcal{K}_c(\mathbb{R})$ verifying that $A = B + C$. The existence of C is subject to the fulfilment of the expression $\text{spr } B \leq \text{spr } A$.

Given the intervals $A, B \in \mathcal{K}_c(\mathbb{R})$, the metric $d_\tau(A, B) = ((1 - \tau)((\text{mid } A - \text{mid } B)^2 + \tau(\text{spr } A - \text{spr } B)^2))^{\frac{1}{2}}$, for an arbitrary $\tau \in (0, 1)$, is the L_2 -type distance to be considered. d_τ is based on the metric d_θ defined in [11].

Given a probability space (Ω, \mathcal{A}, P) the mapping $\mathbf{x} : \Omega \rightarrow \mathcal{K}_c(\mathbb{R})$ is a random interval iff it is a measurable Borel mapping. The moments to be considered are the classical Aumann expected value for intervals; the variance defined following the usual Fréchet variance [8] associated with the Aumann expectation in the interval space $(\mathcal{K}_c(\mathbb{R}), d_\tau)$; and the covariance defined in terms of mids and spreads as $\sigma_{\mathbf{x}, \mathbf{y}} = (1 - \tau) \sigma_{\text{mid } \mathbf{x}, \text{mid } \mathbf{y}} + \tau \sigma_{\text{spr } \mathbf{x}, \text{spr } \mathbf{y}}$.

3 The multiple linear regression model

Let \mathbf{y} be a response random interval and let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$ be k explanatory random intervals. The model is formalized in a matrix notation as follows:

$$\mathbf{y} = X^{Bl} B + \varepsilon, \quad (1)$$

where $B = (b_1|b_2|b_3|b_4)^t \in \mathbb{R}^{4k \times 1}$ with $b_i \in \mathbb{R}^k$ ($i \in \{1, 2, 3, 4\}$), $X^{Bl} = (\mathbf{x}^M|\mathbf{x}^S|\mathbf{x}^C|\mathbf{x}^R) \in \mathcal{K}_c(\mathbb{R})^{1 \times 4k}$ where the elements are defined as $\mathbf{x}^M = \text{mid } x^t [1 \pm 0]$, $\mathbf{x}^S = \text{spr } x^t [0 \pm 1]$, $\mathbf{x}^C = \text{mid } x^t [0 \pm 1]$ and $\mathbf{x}^R = \text{spr } x^t [1 \pm 0]$, considering the canonical decomposition of the regressors.

$\text{mid } x = (\text{mid } \mathbf{x}_1, \text{mid } \mathbf{x}_2, \dots, \text{mid } \mathbf{x}_k)^t \in \mathbb{R}^k$ (analogously $\text{spr } x$) and ε is a random interval-valued error such that $E(\varepsilon|x) = \Delta \in \mathcal{K}_c(\mathbb{R})$.

The following separate linear relationships for the *mid* and *spr* components of the intervals are derived from (1):

$$\text{mid } \mathbf{y} = \text{mid } x^t b_1 + \text{spr } x^t b_4 + \text{mid } \varepsilon, \quad (2a)$$

$$\text{spr } \mathbf{y} = \text{spr } x^t |b_2| + |\text{mid } x^t| |b_3| + \text{spr } \varepsilon. \quad (2b)$$

Thus, the flexibility of the model arises from the possibility of considering all the information provided by $\text{mid } x$ and $\text{spr } x$ to model $\text{mid } \mathbf{y}$ and $\text{spr } \mathbf{y}$, as follows from (2a) and (2b). This represents an improvement with respect to previous models that merely addressed the relationship between the mids of the variables or between the spreads but never any cross-relationship (mid-spr).

Nevertheless, the inclusion of more coefficients entails an increase in the dimensionality of the estimation process. Some of these coefficients could be zero as not all the new introduced variables will contribute. Therefore it is proposed to estimate (1) by least-squares and by Lasso and compare the advantages and disadvantages that each estimation process provides.

3.1 The Least-Squares estimation

Given $\{(\mathbf{y}_j, \mathbf{x}_{i,j}) : i = 1, \dots, k, j = 1, \dots, n\}$ a simple random sample of intervals obtained from $(\mathbf{y}, \mathbf{x}_1, \dots, \mathbf{x}_k)$ in (1) the estimated model is

$$\hat{\mathbf{y}} = X^{ebl} \hat{B} + \hat{\varepsilon} \quad (3)$$

where $y = (\mathbf{y}_1, \dots, \mathbf{y}_n)^t$, $X^{ebl} = (X^M|X^S|X^C|X^R) \in \mathcal{K}_c(\mathbb{R})^{n \times 4k}$, $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^t$ is such that $E(\varepsilon|x) = 1^n \Delta$ and B as in (1). X^M is the $(n \times k)$ -interval-valued matrix such that $(X^M)_{j,i} = \text{mid } \mathbf{x}_{i,j} [1 \pm 0]$ (analogously X^S, X^C and X^R). Given an arbitrary vector of regression coefficients $A \in \mathbb{R}^{4k \times 1}$ and an interval of residuals $C \in \mathcal{K}_c(\mathbb{R})$, the Least-Squares estimation looks for \hat{B} and $\hat{\Delta}$ minimizing the distance $d_\tau^2(y, X^{ebl} A + 1^n C)$. $\hat{\Delta}$ can be obtained separately and firstly by the expression $\hat{\Delta} = \overline{\mathbf{y}} -_H \overline{X^{ebl} \hat{B}}$.

Recalling that, by definition, $X^S = -X^S$ (and analogously $X^C = -X^C$) the estimation process of the coefficients b_2 and b_3 accompanying these variables can be simplified by searching only for non-negative estimates. By contrast,

coefficients b_1 and b_4 are not affected by any kind of restrictions so they can be estimated directly by OLS. Moreover, it has to be assured the existence of the residuals defined as the Hukuhara differences $\varepsilon = \mathbf{y} -_H X^{ebl} B$. For this purpose the minimization problem ends up to be the following constrained quadratic problem:

$$\min_{A_m \in \mathbb{R}^{2k}, A_s \in \Gamma} (1 - \tau) \|v_m - F_m A_m\|^2 + \tau \|v_s - F_s A_s\|^2 \quad (4)$$

$$\Gamma = \{(a_2, a_3) \in [0, \infty)^k \times [0, \infty)^k : \text{spr } X a_2 + |\text{mid } X| a_3 \leq \text{spr } y\},$$

being $v_m = \text{mid } y - \overline{\text{mid } y} 1^n$, $v_s = \text{spr } y - \overline{\text{spr } y} 1^n \in \mathbb{R}^n$, $F_m = \text{mid } X^{ebl} - 1^n (\text{mid } X^{ebl})$, $F_s = \text{spr } X^{ebl} - 1^n (\text{spr } X^{ebl}) \in \mathbb{R}^{n \times 2k}$, $A_m = (a_1 | a_4)^t \in \mathbb{R}^{2k \times 1}$ the coefficients related to the midpoints and $A_s = (a_2 | a_3)^t \in \mathbb{R}^{2k \times 1}$ the coefficients related to the spreads, with $a_l \in \mathbb{R}^k$, $l = 1, \dots, 4$.

There are several numerical ways to tackle the resolution of a quadratic problem as (4). Given the shape of the objective function, the minimization process is solved separately over A_m and A_s . Those coefficients related with the mids (A_m) are not affected by constraints and therefore, the OLS estimator can be used directly. Thus $\widehat{A_m} = (F_m^t F_m)^{-1} F_m^t v_m$. However, in order to proceed with the constrained minimization over A_s , Karush-Kuhn-Tucker conditions guarantee the existence of local optima solution, which can be computed with standard numerical tool. Nevertheless, in order to obtain an exact solution and a more handy estimator of A_s , (4) can be equivalent expressed as a *Linear Complementary Problem* with the shape:

$$\omega = M \lambda + q \quad \text{s.t.} \quad \omega, \lambda \geq 0, \quad \omega_j \lambda_j = 0, \quad j = 1, \dots, n+1, \quad (5)$$

with $M = (R Q^{-1} R^t)$ and $q = (-R Q^{-1} c - r)$ (details in [3]). Thereby, once λ is obtained, the expression of the estimator is $\widehat{A_s} = Q^{-1} (R^t \lambda - c)$.

3.2 The Lasso estimation

Lasso, *Least Absolute Shrinkage and Selection Operator*, is a regression method that involves penalizing the sum of the absolute values of the regression coefficients estimates. For this purpose it involves a regularization parameter which affects directly the estimates: the larger the value of this parameter, the more estimates that are shrunk towards zero. This coefficient cannot be estimated statistically, so a cross-validation process is usually applied.

As previously, (4) can be solved separately. On one hand, the classical Lasso method will be used to obtain the estimator of the regression coefficients related to the mids. Then, the problem is expressed as:

$$\frac{1}{2} \|v_m - A_m F_m\|_2^2 + \lambda \sum_{j=1}^{2k} |A_{m_j}|$$

being λ the regularization parameter. There are different programs capable to solve this problem (such as Matlab or R). The `lasso.m` Matlab function is the one used to obtain $\widehat{A_m}$.

On the other hand, for those coefficients related with the spreads a constrained Lasso algorithm has been developed as a modified version of the code proposed by Mark Schmidt (2005) [10] and is available upon request. The problem is given by:

$$\frac{1}{2} \|v_s - A_s F_s\|_2^2 + \lambda \sum_{j=1}^{2k} |A_{s_j}| \quad \text{s.t. } RA_s \geq r.$$

The most usual elections of λ are the value than minimizes the Cross-Validation Mean Square Error (λ_{MSE}) and the value that provides a simpler or more parsimonious model with respect to λ_{MSE} (in terms of more zero coefficients) but at the same time with one-standard-error (λ_{1SE}).

4 Giordani's Lasso estimation

The so-called *Lasso-based Interval-valued Regression (Lasso-IR)* proposed by Giordani in [4] is another Lasso method to deal with a multiple linear regression model for interval data. However, the later regression model is not formalized following the interval arithmetic and can end up with an ill-defined estimated model. Keeping the same notation as in (2b), it requires the non-negativity of b_2 and b_3 but does not test if the Hukuhara's difference $\varepsilon = \mathbf{y} -_H X^{ebl} B$ exists. The optimization problem can be written (analogously to (4)) as:

$$\min_{A_m, A_s} (1 - \tau) \|v_m - F_m A_m\|^2 + \tau \|v_s - F_s(A_m + A_a)\|^2 \quad (6)$$

$$F_s(A_m + A_a) \geq 0, \sum_{j=0}^p |A_{a_j}| \leq t$$

The coefficients related to the spreads (A_s) are the ones for the mids (A_m) plus a vector of additive coefficients (A_a) showing the distance that they are allowed to differ from A_m . In this case (6) has been expressed as a constrained quadratic problem, where there is a one-to-one correspondence between λ and t . The value of t that minimizes the cross-validation mean square error is the one considered. In order to solve the problem a stepwise algorithm based on [6] is proposed.

Another important difference, which entails less flexibility in the model, is the limitation of being able to study separately the relationships between the mids and the relationship between the spreads of the intervals but never any cross-relationship.

Remark 1. There is a particular case of model (1), the so-called Model M addressed in [2], which is formalized in the interval framework but has the same lack of flexibility as (6). In this case b_3 and $b_4 = (0, \dots, 0)$, so the model has the shape:

$$\mathbf{y} = b_1 \mathbf{x}^M + b_2 \mathbf{x}^S + \varepsilon. \quad (7)$$

5 A real-life illustrative example

The following example contains the information of a sample of 59 patients (from a population of 3000) hospitalized in the Hospital Valle del Nalón in Asturias, Spain. The variables to be considered are the ranges of fluctuation of the diastolic blood pressure over the day (\mathbf{y}), the pulse rate (\mathbf{x}_1) and the systolic blood pressure (\mathbf{x}_2). The dataset can be found in [2] and [5].

In order to make possible the comparison between the estimator proposed in Sect. 4 and those ones introduced in Subsect. 3.1 and Subsect. 3.2, the example will be developed for the simpler model explained in Remark 1.

Given the displayed model in (7), $\mathbf{y} = b_1\mathbf{x}_1^M + b_2\mathbf{x}_2^M + b_3\mathbf{x}_1^S + b_4\mathbf{x}_2^S + \varepsilon$, the estimates of the regression coefficients are summarized in Table 1:

Table 1. Estimates of the regression coefficients for the three estimators: LS, Lasso (for the two more representatives values of λ) and Lasso-IR (for a fixed value of $t=0.10$ prefixed by the author). The last column contains the MSE of the models mimicking its definition in the classical framework.

	\hat{b}_1	\hat{b}_2	\hat{b}_3	\hat{b}_4	MSE
<i>LS – estimation (Sect. 3.1)</i>	0.4497	0.0517	0.2588	0.1685	68.2072
<i>Lasso – estimation (Sect. 3.2)</i>	0.4202	0.0020	0.3379	0.2189	68.8477
λ_{MSE}	(0.6094)		(0.0259)		
<i>Lasso – estimation (Sect. 3.2)</i>	0.2749	0	0.0815	0	76.9950
λ_{1SE}	(3.2521)		(1.8736)		
<i>Lasso – IR (Sect. 4)</i>	0.5038	0.1261	0.4847	0.3605	71.2418

In view of the results in Table 1, those coefficients which take small values with the LS-estimation (\hat{b}_2 and \hat{b}_4) are shrunk towards zero with the most preferable Lasso estimation (for λ_{1SE}). However, this entails a significant increase of the MSE. In the case of using our Lasso-estimator for λ_{MSE} , the MSE is smaller but it does not provide a parsimonious model, being therefore its usefulness questionable. The estimator proposed in [4] reaches a high value of MSE (worse in comparison with the lasso for λ_{MSE}) and does not end up with an easy-to-interpret model.

6 Conclusions

On one hand, a recently studied regression model for interval data, allowing to study all the cross-relationships between the mids and spreads of the interval-valued variables involved, is considered. This flexibility derives into an increase of the dimensionality of the model. Therefore a Lasso estimation seems appropriate to tackle this problem by setting some of these coefficients to zero. Nonetheless, a comparison study gathering the double estimation process conducted (first by Least-Squares and after by Lasso) is provided.

On the other hand, it is considered the Lasso-based interval-valued regression model (Lasso-IR) proposed in [4]. This model is not constrained to guarantee the existence of the residuals so it can provide misleading estimations. Moreover, it has a lack of flexibility as it solely tackles the relationships of type mid-mid and spr-spr but no cross-relationship mid-spr.

A real-life example illustrating the difference between the estimators in terms of MSE and simplicity has been conducted.

Acknowledgments. Financial support from the Spanish Ministerio de Ciencia e Innovación (MICINN) through Ayuda Puente SV-PA-13-ECOEMP-66 and Acción Integrada PRI-AIBDE-2011-1197, is kindly acknowledged.

References

1. Abdallah, F., Gning, A., Bonnifait, P.: Adapting particle filter on interval data for dynamic state estimation. In: ICASSP, pp. 1153-1156, (2007)
2. Blanco-Fernández, A., Corral, N., González-Rodríguez, G.: Estimation of a flexible simple linear model for interval data based on set arithmetic. *Comput. Stat. Data Anal.* 55, 2568-2578 (2011)
3. Blanco-Fernández, A., García-Bárcana, M., Colubi, A., Kontoghiorghe, E.J.: Multiple set arithmetic-based linear regression models for interval-valued variables (submitted)
4. Giordani, P.: Linear regression analysis for interval-valued data based on the Lasso technique. In: 58th Session of the International Statistical Institute, pp. 5576-5581, Ireland (2011)
5. González-Rodríguez, G., Blanco, A., Corral, N., Colubi, A.: Least squares estimation of linear regression models for convex compact random sets. *Adv. Data Anal. Classif.* 1, 67-81 (2007)
6. Lawson, C.L., Hanson, R.J.: Solving Least Squares Problems. In: *Classics in Applied Mathematics*, vol. 15. SIAM, Philadelphia (1995)
7. Lima Neto, E.A., de Carvalho, F.A.T.: Constrained linear regression models for symbolic interval-valued variables. *Comput. Stat. Data Anal.* 54, 333-347 (2010)
8. Näther, W.: Linear statistical inference for random fuzzy data. *Statistics.* 29(3), 221-240 (1997)
9. Ramos-Guajardo, A.B., González-Rodríguez, G.: Testing the Variability of Interval Data: An Application to Tidal Fluctuation. In: Borgelt, C., Gil, M.A., Sousa, J.M.C., Verleysen, M. (eds.) *Towards Advanced Data Analysis by Combining Soft Computing and Statistics. Studies in Fuzziness and Soft Computing*, vol. 285, pp. 65-74 (2013)
10. Schmidt Mark, <http://www.di.ens.fr/~mschmidt>
11. Trutschnig, W., González-Rodríguez, G., Colubi, A., Gil, M.A.: A new family of metrics for compact, convex (fuzzy) sets based on a generalized concept of mid and spread. *Inf. Sci.* 179(23), 3964-3972 (2009)